

GPU-based Real-time Triggering in the NA62 Experiment

R. Ammendola[†], A. Biagioni^{*}, P. Cretaro^{*}, S. Di Lorenzo[‡], R. Fantechi[‡], M. Fiorini^{**}, O. Frezza^{*}, G. Lamanna^{‡§}, F. Lo Cicero^{*}, A. Lonardo^{*}, M. Martinelli^{*}, I. Neri^{**}, P. S. Paolucci^{*}, E. Pastorelli^{*}, R. Piandani[‡], L. Pontisso[‡], D. Rossetti^{||}, F. Simula^{*}, M. Sozzi^{‡§}, and P. Vicini^{*}

Abstract—Over the last few years the GPGPU (General-Purpose computing on Graphics Processing Units) paradigm represented a remarkable development in the world of computing. Computing for High-Energy Physics is no exception: several works have demonstrated the effectiveness of the integration of GPU-based systems in high level trigger of different experiments. On the other hand the use of GPUs in the low level trigger systems, characterized by stringent real-time constraints, such as tight time budget and high throughput, poses several challenges. In this paper we focus on the low level trigger in the CERN NA62 experiment, investigating the use of real-time computing on GPUs in this synchronous system. Our approach aimed at harvesting the GPU computing power to build in real-time refined physics-related trigger primitives for the RICH detector, as the the knowledge of Cerenkov rings parameters allows to build stringent conditions for data selection at trigger level. Latencies of all components of the trigger chain have been analyzed, pointing out that networking is the most critical one. To keep the latency of data transfer task under control, we devised NaNet, an FPGA-based PCIe Network Interface Card (NIC) with GPUDirect capabilities. For the processing task, we developed specific multiple ring trigger algorithms to leverage the parallel architecture of GPUs and increase the processing throughput to keep up with the high event rate. Results obtained during the first months of 2016 NA62 run are presented and discussed.

I. INTRODUCTION

In High Energy Physics experiments the realtime selection of the most interesting events is of paramount importance because of the collision rates which do not give the possibility to save all the data for offline analysis. For this purpose, different trigger levels are usually used to carefully choose the most meaningful events. The low level ones require low and (almost) deterministic latency and their standard implementation is on dedicated hardware (ASICs or FPGAs). Our approach aims at exploiting the Graphic Processing Units (GPUs) computing power, in order to build refined physics-related trigger primitives, such as energy or direction of the final state particles in the detectors, and therefore leading to a net improvement of trigger conditions and data handling. GPUs architectures are massively parallel, being designed to optimize computing throughput but with no particular attention to their usage in real-time contexts, such as

the online low level triggers. While execution times are rather stable on these devices, also I/O tasks have to guarantee real-time features along the data stream path, from detectors to GPU memories. The NaNet project arises with the goal of designing a low-latency and high-throughput data transport mechanism for systems based on CPU/GPUs. The GPU-based L0 trigger using the NaNet board is currently integrated in the experimental setup of the RICH Čerenkov detector of the NA62 experiment in order to reconstruct the ring-shaped hit patterns. We report and discuss results obtained with this system along with the algorithms that will be implemented.

II. NANET ARCHITECTURE

The design of a low-latency, high-throughput data transport mechanism for real-time systems is mandatory in order to bridge the front-end electronics and the software trigger computing nodes [1] of High Energy Physics Experiments. NaNet, being an FPGA-based NIC, natively supports a variety of link technologies allowing for a straightforward integration in different experimental setups. Its key characteristic is the management of custom and standard network protocols in hardware, in order to avoid OS jitter effects and guarantee a deterministic behaviour of communication latency while achieving maximum capability of the adopted channel. Furthermore, NaNet integrates a processing stage which is able to reorganize data coming from detectors on the fly, in order to improve the efficiency of applications running on computing nodes. On a per-experiment basis, different solutions can be implemented: data decompression, reformatting, merging of event fragments.

Finally, data transfers to or from application memory are directly managed avoiding bounce buffers. NaNet accomplishes this zero-copy networking by means of a hardware implemented memory copy engine that follows the RDMA paradigm for both CPU and GPU — this latter supporting the GPUDirect V2/RDMA by NVIDIA to minimize the I/O latency in communicating with GPU accelerators.

On the host side, a dedicated Linux kernel driver offers its services to an application level library, which provides the user with a series of functions to: open/close the NaNet device; register and de-register circular lists of persistent data receiving buffers (CLOPs) in GPU and/or host memory; manage software events generated when a receiving CLOP buffer is full (or when a configurable timeout is reached) and received data are ready to be processed.

^{*}INFN Sezione di Roma, Italy.

[†]INFN Sezione di Tor Vergata, Italy.

[‡]INFN Sezione di Pisa, Italy.

[§]CERN, Switzerland.

^{||}Università di Roma Sapienza, Dipartimento di Fisica, Italy.

^{||}NVIDIA Corporation, U.S.A.

^{**}INFN Sezione di Ferrara, Italy.

NaNet-1 was developed in order to verify the feasibility of the project; it is a PCIe Gen2 x8 network interface card featuring GPUDirect RDMA over GbE.

NaNet-10 is a PCIe Gen2 x8 network adapter implemented on the Terasic DE5-net board equipped with an Altera Stratix V FPGA featuring four SFP+ cages [2].

Both implementations use UDP as transport protocol.

In Fig. 1, NaNet-10 and NaNet-1 latencies are compared within UDP datagram size range; NaNet-10 guarantees sub- μ s hardware latency for buffers up to ~ 1 kByte in GPU/CPU and it reaches its 10 Gbps bandwidth peak capability already at ~ 1 kByte size (Fig. 2).

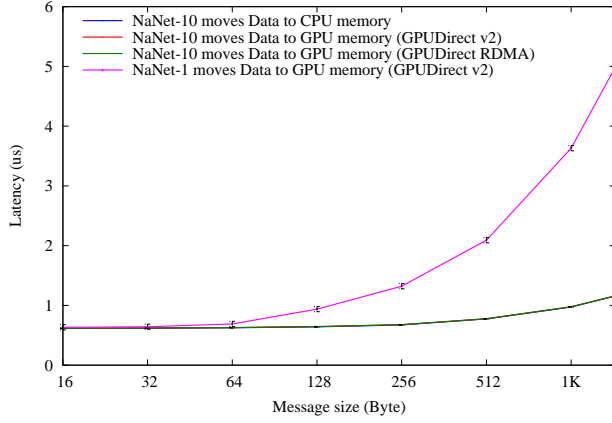


Fig. 1. NaNet-10 vs. NaNet-1 hardware latency. NaNet-10 curves are completely overlapping at this scale.

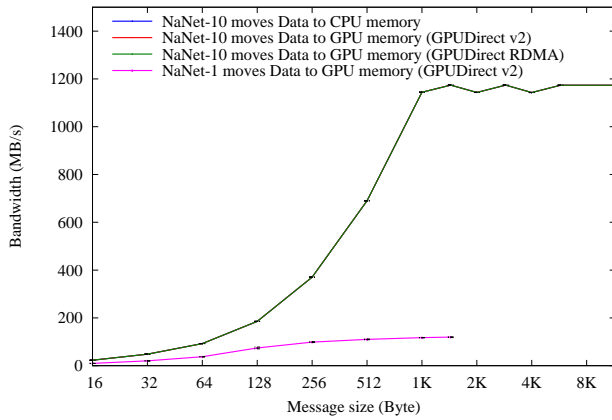


Fig. 2. NaNet-10 vs. NaNet-1 bandwidth. NaNet-10 curves are completely overlapping at this scale.

III. NA62 EXPERIMENT

The NA62 experiment at CERN [3] aims at measuring the Branching Ratio of the ultra-rare decay of the charged Kaon into a pion and a neutrino-antineutrino pair. The NA62 goal is to collect ~ 100 events with a signal to background ratio 10:1, using a novel technique with a high-energy (75 GeV)

unseparated hadron beam decaying in flight. In order to manage the high-rate data stream due to a ~ 10 MHz rate of particle decays illuminating the detectors, a set of trigger levels will have to reduce this rate by three orders of magnitude. The entire trigger chain works on the main digitized data stream [4].

The Low-level trigger, implemented in hardware by means of FPGAs on the readout boards, reduces the data stream by a factor 10 to meet the maximum design rate for event readout of 1 MHz. The upper trigger levels (L1 and L2) are software-implemented on a commodity PC farm for further reconstruction and event building.

In the standard implementation, the FPGAs on the readout boards compute simple trigger primitives on the fly, such as hit multiplicities and rough hit patterns, which are then time-stamped and sent to a central processor for matching and trigger decision. Thus the maximum latency allowed for the synchronous L0 trigger is related to the maximum data storage time available on the data acquisition boards. For NA62 this value is up to 1 ms, in principle allowing use of more compute demanding implementations at this level, *i.e.* the GPUs.

IV. IMPLEMENTATION OF THE GPU-BASED LOW-LEVEL TRIGGER

As a first example of GPU application in the NA62 trigger system we studied the possibility to reconstruct rings in the RICH. This detector identifies pions and muons with momentum in the range between 15 GeV/c and 35 GeV/c. Čerenkov light is reflected by a composite mirror with a focal length of 17 m focused onto two separated spots equipped with ~ 1000 photomultipliers (PM) each. Data communication between the readout boards (TEL62) and the L0 trigger processor happens over multiple GbE links using UDP streams. The final system consists of 4 GbE links to move primitives data from the readout boards to the GPU_L0TP (see Fig. 3). The overall time budget for the low level trigger comprising both communication and computation tasks is of 1 ms, so a deterministic response latency from GPU_L0TP is a strict requirement. Refined primitives coming from the GPU-based calculation will be then sent to the central L0 processor, where the trigger decision is made taking in account informations from other detectors.

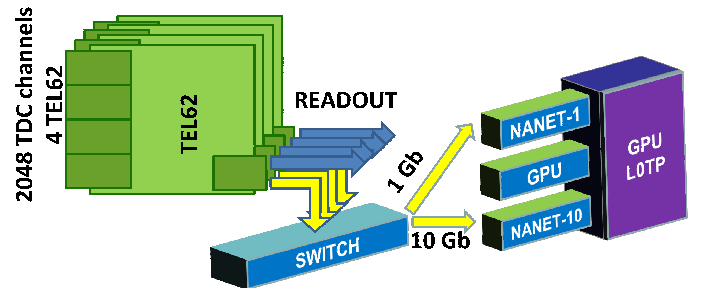


Fig. 3. Pictorial view of GPU-based Trigger.

A. Multiple ring events reconstruction on GPU

Taking the parameters of Čerenkov rings into account could be very useful in order to build stringent conditions for data selection at trigger level. This implies that circles have to be reconstructed using the coordinates of activated PMs.

We take in consideration two multi-rings pattern recognition algorithms based only on geometrical considerations (no other information is available at this level) and particularly suitable for exploiting the intrinsic parallel architecture of GPUs.

1) *Histogram algorithm*: The procedure involves dividing the XY plane into a grid and creating a histogram whose bins contain distances from the gridpoints and the hits of the physics event. Distance bins whose contents exceed a threshold value let identify the rings. In order to limit the use of resources, it is possible to proceed in two steps, starting the histogram procedure with a 8x8 grid and calculating now distances from such squares. Afterwards, to refine their positions, the calculation is repeated with a grid 2x2 only for the candidate centers selected according to the threshold in the previous step.

2) *Almagest algorithm*: The Ptolemy's Theorem states that when four vertices of a quadrilateral (ABCD) lie on a common circle, it is possible to relate four sides and two diagonals: $|AC| \times |BD| = |AB| \times |CD| + |BC| \times |AD|$. By using this formula it is possible to implement a pattern recognition algorithm for multi-rings which exposes different level of parallelism, resulting well-suited for GPUs architecture and fast in its execution. This is crucial either to directly reconstruct the rings or to choose different algorithms according to the number of circles. The large number of possible combinations of four vertices, given a maximum of 64 points for physics event, can be a limitation to this approach. To greatly reduce the number of tests, one possibility is to choose few triplets — *i.e.* a set of three hits assumed to belong to the same ring — trying to maximize the probability that all their points belong to the same ring and iterating through all the remaining hits to search for the ones satisfying the aforementioned formula [5]. The parallel implementation of this algorithm yields many triplets and events being processed at the same time. Some results are shown in Fig. 4, where the computing time is further sped up by greatly reducing accesses to GPU shared memory, mainly using threads private registers through CUDA intra-warp veto and shuffle instructions, so that multi-rings events are processed with a latency of 0.5 us per event. Once the number of rings and points belonging to them have been found, it is possible to apply *e.g.* Crawford's method [6] to obtain centre coordinates and radii with better spatial resolution.

V. RESULTS

In 2015 the GPU-based trigger at CERN comprises 2 TEL62 readout boards connected to a HP2920 switch and a NaNet-1 board with a TTC HSMC daughtercard plugged into a SuperMicro server consisting of a X9DRG-QF dual socket motherboard — Intel C602 Patsburg chipset — populated with Intel Xeon E5-2620 @2.00 GHz CPUs (*i.e.* Ivy Bridge micro-architecture), 32 GB of DDR3 memory and a Kepler-class NVIDIA K20c GPU.

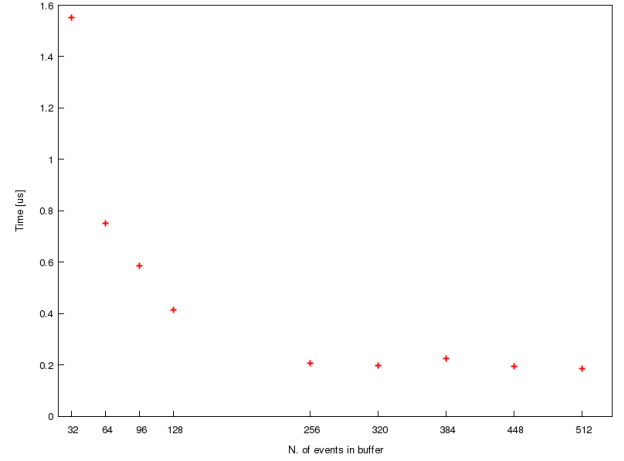


Fig. 4. Almagest algorithm performances, time for single event. Test performed on K20c NVIDIA GPU.

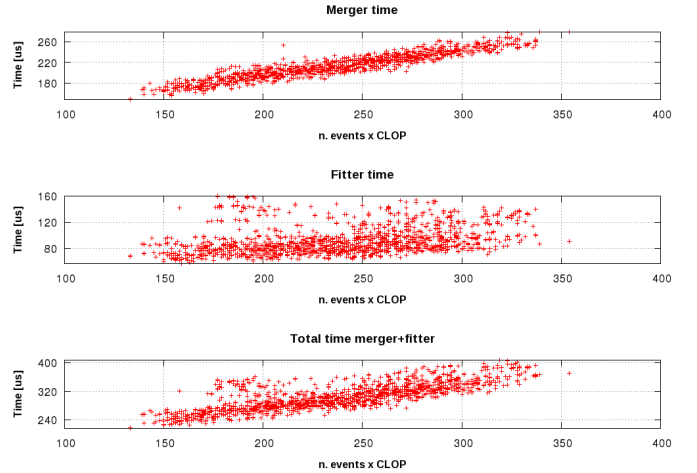


Fig. 5. Multi-ring reconstruction of events performed on K20c NVIDIA GPU.

Such a system allows for testing of the whole chain: the data events move towards the GPU-based trigger through NaNet-1 by means of the GPUDirect RDMA interface. Data arriving within a configurable time frame are gathered and then organized in a Circular List Of Persistent buffers (CLOP) in the GPU memory. Buffer number and size are tunable in order to optimize computing and communication. This time frame must obviously be shorter or equal on average to how long the GPU takes for multi-ring reconstruction, to be sure that buffers are not overwritten by incoming events before they are consumed by the GPU. Events coming from different TEL62 need to be merged in the GPU memory before the launch of the ring reconstruction kernel. Each event is timestamped and the ones coming from different readout boards that are in the same time-window are fused in a single event describing the status of PMs in the RICH detector.

Results are reported in Fig. 5. The CLOP size measured as number of received events is on the X-axis and the latencies of different stages are on the Y-axis. The computing kernel implemented the histogram fitter with a single step (*i.e.* using

an 8x8 grid only). Events coming from 2 readout boards, for a gathering time of 400 μ s, and parameters like events rate (collected with a beam intensity of 4×10^{11} protons per spill), a CLOP's size of 8KB, time frame was chosen so that we could test the online behaviour of the trigger chain.

Because the merge operation doesn't expose much parallelism, requiring instead synchronization and serialization, this is an ill-suited problem to the GPU architecture. In operative conditions, the merging time only would exceed the time frame. The high latency of the merger task when performed on a GPU strongly suggests to offload such duties to a hardware implementation.

VI. CONCLUSIONS AND FUTURE WORK

ACKNOWLEDGMENT

S. Di Lorenzo, R. Fantechi, M. Fiorini, I. Neri, R. Piandani, L. Pontisso, M. Sozzi thank the GAP project, partially supported by MIUR under grant RBFR12JF2Z "Futuro in ricerca 2012".

REFERENCES

- [1] A. Lonardo *et al.*, "NaNet: a Configurable NIC Bridging the Gap Between HPC and Real-time HEP GPU Computing," *Journal of Instrumentation*, vol. 10, no. 04, p. C04011, 2015. [Online]. Available: <http://stacks.iop.org/1748-0221/10/i=04/a=C04011>
- [2] R. Ammendola, A. Biagioni, M. Fiorini, O. Frezza, A. Lonardo, G. Lamanna, F. Lo Cicero, M. Martinelli, I. Neri, P. Paolucci, E. Pastorelli, L. Pontisso, D. Rossetti, F. Simula, M. Sozzi, L. Tosoratto, and P. Vicini, "Nanet-10: a 10gbe network interface card for the gpu-based low-level trigger of the na62 rich detector," *Journal of Instrumentation*, vol. 11, no. 03, p. C03030, 2016. [Online]. Available: <http://stacks.iop.org/1748-0221/11/i=03/a=C03030>
- [3] G. Lamanna, "The NA62 experiment at CERN," *Journal of Physics: Conference Series*, vol. 335, no. 1, p. 012071, 2011. [Online]. Available: <http://stacks.iop.org/1742-6596/335/i=1/a=012071>
- [4] C. Avanzini *et al.*, "The trigger and DAQ system for the NA62 experiment," *Nucl. Instrum. Methods Phys. Res., A*, vol. 623, pp. 543–545, 2010.
- [5] G. Lamanna, "Almagest, a new trackless ring finding algorithm," *Nuclear Instruments and Methods in Physics Research Section A: Accelerators, Spectrometers, Detectors and Associated Equipment*, vol. 766, pp. 241 – 244, 2014, {RICH2013} Proceedings of the Eighth International Workshop on Ring Imaging Cherenkov Detectors Shonan, Kanagawa, Japan, December 2-6, 2013. [Online]. Available: <http://www.sciencedirect.com/science/article/pii/S0168900214006135>
- [6] J. Crawford, "A non-iterative method for fitting circular arcs to measured points," *Nuclear Instruments and Methods in Physics Research*, vol. 211, no. 1, pp. 223 – 225, 1983. [Online]. Available: <http://www.sciencedirect.com/science/article/pii/0167508783905756>